

White Paper Report

Report ID: 109287

Application Number: HD-51636-13

Project Director: Susan Weiss (sweiss@jhu.edu)

Institution: Johns Hopkins University

Reporting Period: 5/1/2013-10/31/2015

Report Due: 1/31/2016

Date Submitted: 2/1/2016

NEH WHITE PAPER

White Paper

Report ID: 113804

Application Number: HD-51636-13

Project Director: Susan Weiss (sweiss@jhu.edu) Institution: Johns Hopkins University

Report ID: 109289

Application Number: HD-51636-13

Project Director: Susan Weiss (sweiss@jhu.edu) Institution: Johns Hopkins University

Reporting Period: 11/1/2014-1/31/2016

Report Due: 1/31/2016

Date Submitted: 1//31/2016

Digital Prosopography for Renaissance Musicians: Discovery of Social and Professional
Networks

Susan Weiss and Ichiro Fujinaga

Johns Hopkins University 31 January 2016

a. Project Activities

- **Software Development**

In October 2015, following two semesters of work, our Johns Hopkins graduate students completed creating a web-based work-flow environment with an easy-to-use flexible editor of RDF data (quads) and a simple interface to SPARQL. At that point, because of a number of challenges (code spread over a number of different repositories, a number of coders, some of them not fluent in English who were documenting in Chinese), and, as a result of the bugs and flaws created by these complexities, the principals sought the professional assistance of a former graduate student at the University of Maryland, Max Morawski.

What we discovered was that the dependencies had to be figured out in order to rebuild the project from scratch. Max took on the following:

1. A web interface for uploading text files to be processed.
2. NLP programs in Java (OpenIE, nlptools, etc.) that extract relationships (subject-predicate-object) from the text.
3. Import the resulting triples into the RDF Editor.

He then moved all code in one repository and began restructuring to bring it more in line with standard practices. He also made it possible to work on the entire project with little effort using an Eclipse project and Maven. In the third and fourth week of his one month on the project, he was able to make some critical repairs such as fixing security flaws, and made some cosmetic and developer changes such as the removal of extraneous files (over a dozen), and the correction of typos. He also added new features including passing the RDF triples to the RDF editor after inspection and set Open IE libraries to preload. As a result of his work, upload.html now does the actual RDF parsing and file uploading is

simpler and the text of the files is displayed. The website is hosted on Amazon Web Service including a page for uploading and deleting files, a code for parsing text into RDF triples using Open IE (not attached to the website) and an RDF event editor

- **Project Activities: Knowledge Mobilization, Publicizing**

All throughout the process our findings have been available on (<http://www.humanhistoryproject.ca>). In addition, Fujinaga and Weiss gave a poster presentation entitled: “Digital Prosopography of Renaissance Musicians: Discovery of social and professional network” at the Annual Meeting of the American Musicological Society / Society for Music Theory in Milwaukee on 6 November 2014. We had excellent feedback. We have broadened our base of musicologists in an effort to begin to answer queries. We were also awarded a no-cost extension on the project through to next October 2015 in an effort to finish programming, mine data and continue software development. Our team is looking at ways of solving the problem of teaching the machine to recognize predicates. We are currently experimenting with NLP using reVerb. Our presentation was well-received at the RSA meeting in Berlin in March 2015 where we had queries from musicologists and from scholars in other disciplines. We presented our latest results at the International Association of Music Libraries, Archives and Documentation Centres (IAML) / International Musicological Society (IMS) Congress in New York City on 22 June 2015 at the Juilliard School. Several questions were raised including one about foreign languages. Coincidentally, at that same meeting, we met researchers from abroad who expressed interest in partnering with us in future symposia. Another of our abstracts was submitted and accepted for the Conference on Interdisciplinary Musicology “Imagination in Music, Shanghai, China, 27–29 November.

(see Appendix). In December, Professor Fujinaga hosted a workshop— The second international workshop on Human History Project: Natural language processing and big data— at the Schulich School of Music at McGill University in Montréal (http://www.cirmmt.org/activities/workshops/research/human_history_project_2/hpp).

This workshop was organized by Research Axis 2 of CIRMMT (Centre for Interdisciplinary Research on Music Media and Technology) and was the second workshop on a large project called Human History Project, which aims to build a distributed international database of documented human history using Natural Language Processing tools and Linked Open Data to model historical data. The speakers were:

Jason Boyd, Department of English, Ryerson University
Serge ter Braake, Faculties of Humanities, University of Amsterdam
Ichiro Fujinaga, CIRMMT, Schulich School of Music, McGill University
Sergio Oramas, Music Technology Group, Universitat Pompeu Fabra
Matthew Milner, Department of History, McGill University
Susan Forscher Weiss, Peabody Conservatory of Music, Johns Hopkins University.
Following the workshop the speakers met to map out next steps.

In addition to the conference presentations, we gave seven invited talks regarding this project:

University of Kyoto, Kyoto, Japan. 28 January 2016. *Research Program of the Distributed Digital Music Archives and Libraries Laboratory*.
Japan Advanced Institute of Science and Technology (JAIST), Kanazawa, Japan. 21 January 2016. *Human History Project*.
School of Information Studies Research Seminar, McGill University, Montreal. 13 November 2015. *The Research Program of the Distributed Digital Music Archives and Libraries Laboratory*.
CCS Colloquium, University of Virginia. 2 October 2015. *10 Years of Innovation in Music Digitization & Dissemination*.
Oxford eResearch Centre Seminar. 28 July 2015. *Research Program of the Distributed Digital Music Archives and Libraries Laboratory*.
McGill Digital Humanities Work in Progress Talk, McGill University, Montreal. 9 February 2014. *Human History Project: Storing documented human history into a database*.
École Polytechnique de Montréal, Montreal. 19 November 2013. *Digital Prosopography of Renaissance Musicians*.

b. Accomplishments

- Project web site at humanhistoryproject.ca
- Conducted three major experiments
- Four conference presentations: RSA (2014), AMS/SMT (2014), RSA (2015), and IAML/IMS (2015)
- Two workshops held at the Centre for Interdisciplinary Research on Music Media and Technology (CIRMMT) in Montreal on 21 June 2013 and 12 December 2015.
- Seven invited talks (three in Montreal, one in the USA, and two in Japan).
- Training of several graduate students

c. Audiences

The audiences for this project are musicologists, music librarians, Renaissance scholars, digital humanists, and historians in general. Students in courses across the humanities will benefit from the ability to ask queries across disciplines.

d. Evaluations

Evaluations of our system have been accomplished throughout the project by software testing, software quality evaluations, and several formal experiments using the system.

e. Continuation of the Project

The continuation of the project will be realized through applications to several research grants, publicizing the project, gathering other scholars to the project, and encouraging our graduate students to work on the project.

In November 2015, with Fujinaga as the principal investigator, we have applied to the Quebec government (FRQSC) for a 4-year research grant (\$640,000CAN) entitled: “Music Information, Research, and Infrastructure (MIRAI)”. The study of prosopography of musicians is one of the five major themes of this infrastructure grant.

We have submitted a paper for the Sixteenth Century Society and Conference 2016: ‘Cultural networks in the Renaissance: methodological challenges’ to be held in Bruges, Belgium, 18–20 August 2016. We also plan to submit another paper for the 2nd International Workshop on Computational History and Data-Driven Humanities, to be held in Dublin, 25 May 2016.

f. Long Term Impact

The long-term impact of this project will be realized through continual development of the system and participation of other scholars to test and use the system. As more data is processed by the system, its utility should prove invaluable, as people will be able to search the database for historical events and discover network of people’s relationships in the past.

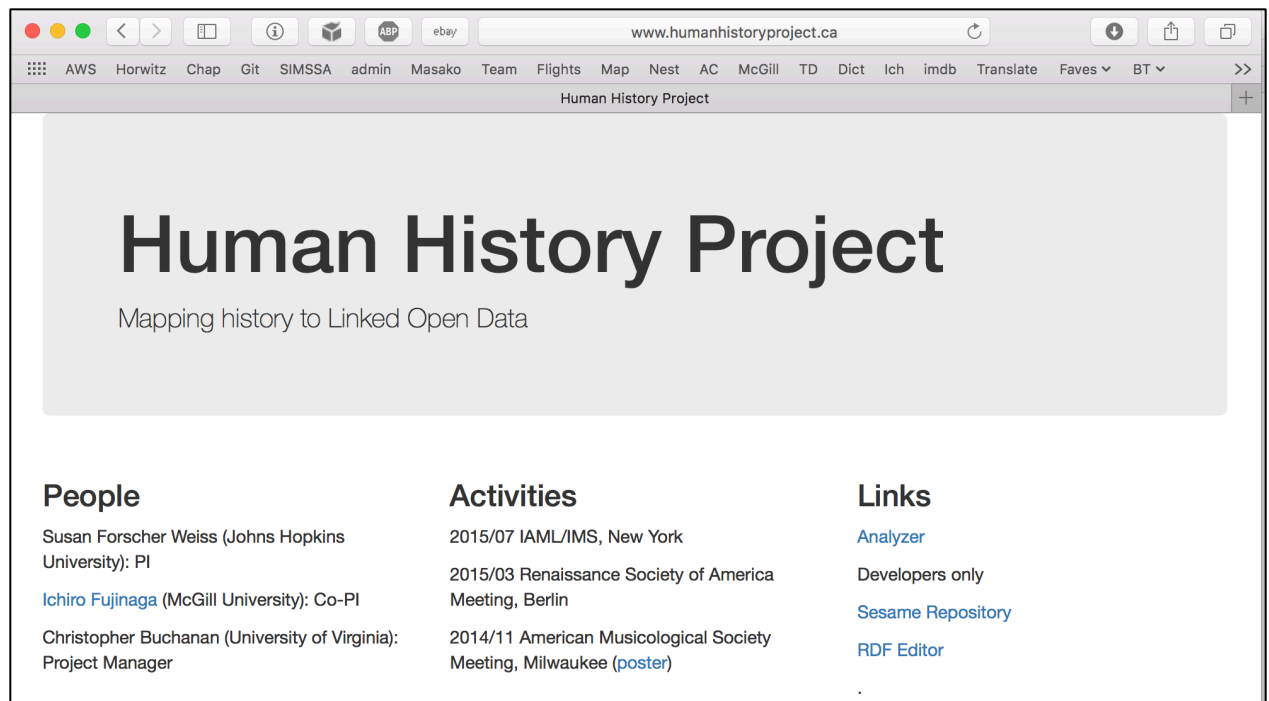
g. Grant Products

- Project web site at humanhistoryproject.ca
- Conference presentations: abstracts, a poster, and several Power Point presentations

3. Appendix

- A screenshot of the website: The Home page
- A screenshot of the website: The Analyzer page
- A screenshot of the website: The RDF Editor page
- Our abstract in CIM Conference Proceedings: 1/2
- Our abstract in CIM Conference Proceedings: 2/2
- Our poster at AMS/SMT 2014
- CIRMMT Workshop #1 The Science and Technology of Music website (2013)
- CIRMMT Workshop #2 The Science and Technology of Music website (2013)
- Slides from our Presentation 1/4
- Slides from our Presentation 2/4
- Slides from our Presentation 3/4
- Slides from our Presentation 4/4

Appendix



A screenshot of the website: The Home page

RDF Web Editor Test Version

Choose File

no file selected

Export

downloadjson

Prefix	Namespace	Select
rdf:	http://www.w3.org/1999/02/22-rdf-syntax-ns#	
rdfs:	http://www.w3.org/2000/01/rdf-schema#	
foaf:	http://xmlns.com/foaf/0.1/	
owl:	http://www.w3.org/2002/07/owl#	
xsd:	http://www.w3.org/2001/XMLSchema#	
Dbpedia-Owl:	http://dbpedia.org/ontology/	
category:	http://dbpedia.org/resource/Category:	
dbpedia:	http://dbpedia.org/resource/	
dbpprop:	http://dbpedia.org/property/	
yago:	http://dbpedia.org/class/yago/	
dcterms:	http://purl.org/dc/terms/	

Add

Delete

Prefix/URL Switch:

Prefix

Full URL

Index	Subject	Predicate	Object	Graph Name	Select
1	Heinrich Glarean	was	a Swiss music theorist, poet and humanist		<input type="checkbox"/>
2	Heinrich Glarean	was born	in Mollis		<input type="checkbox"/>
3	Heinrich Glarean	died	in Freiburg		<input type="checkbox"/>
4	Heinrich Glarean	enrolled	in the University of Cologne		<input type="checkbox"/>
5	Heinrich Glarean	studied	theology, philosophy, and mathematics as well as music		<input type="checkbox"/>
6	It	was	there		<input type="checkbox"/>
7	Heinrich Glarean	wrote	a famous poem		<input type="checkbox"/>
8	Heinrich Glarean	met	Erasmus and the two humanists		<input type="checkbox"/>
9	Heinrich Glarean	became	lifelong friends		<input type="checkbox"/>

Add

Delete

Save to Repo

A screenshot of the website: The RDF Editor page

Imagining the musical past: Creating a digital prosopography of Renaissance musicians

Susan Forscher Weiss (Peabody Conservatory of Music, The Johns Hopkins University)
Ichiro Fujinaga (McGill University)

Background in Music History, German & Romance Languages and Literature
Background in Music Technology

Aims

We imagine the music of the past by attempting to perform in as authentic a manner as possible: by building close models of the musical instruments, by studying drawings and paintings, and by close reading of the contemporary treatises. We can also enhance our image of the past musical world by studying the people and their social or professional network based on multiple sources of information. Such a study is called prosopography. By gathering lots of small bits of information, historians can imagine and make sense of the past (Collingwood 1946; Wiseman 1994).

Main Contribution

Currently, building a database of people's relationships is costly and time consuming because it can only be done manually. This project aims to create, automatically, an economically feasible digital prosopographical database of Renaissance musicians by exploiting the ever-increasing availability of historical documents, recent improvements in optical character recognition (OCR) and natural language processing (NLP) technologies, and the emerging data structure called Linked Open Data.

We are creating a framework that can answer questions not easily answered by Google-like searches or traditional means. For example, which printers in Venice in the 1530s were publishing books of music? Which foreign musicians visited Venice in 1538? Did composer A and composer B live in Venice in 1538? Were there musicians working in Venice from 1535–1540 who performed music by both of these composers? Who were the musical instrument makers there in those years?

The possibility of digitizing prosopography has been discussed since the late 1980's (e.g., Bulst 1989). One of the earliest examples of digital prosopography is the *Prosopography of the Byzantine Empire I* (Bradley & Short 2001), produced on CD-ROM. The largest example of digital prosopography is the *Oxford Dictionary of National Biography*. It is a monumental work, which took over 10 years (1992–2004) to complete at a cost of over £25million (Harrison 2004). More recent digital prosopography projects include the Berkeley Prosopography Services (<http://berkeleyprosopography.org>), the Prosopographie des Chantres de la Renaissance (<http://ricercar.cesr.univ-tours.fr/3-programmes/PCR/>), and the London Lives Project. The latter project, which contains over 3 million names, scanned microfilm sources and entered texts manually (double keyed). The names, occupations, places, and dates were

marked up using “a combination of automated and manual processes.”¹ In our project we are further automating the process by using state-of-the-art OCR and NLP technologies. This should drastically reduce the cost compared to previous projects. Because of errors as a result of imperfections in OCR and NLP technologies we plan to deploy crowd- or expert-sourcing techniques for corrections.

The main types of information we are interested in extracting are named entities (person, place, organization, etc.), events, and relationships between named entities and events. GATE (General Architecture for Text Engineering) is our main NLP tool but we are also working with Stanford NLP software (<http://nlp.stanford.edu/software/>), UIMA (Unstructured Information Management Applications; <http://uima.apache.org>), and REEL (RElationship Extraction Learning framework; <http://reel.cs.columbia.edu>) in an effort to improve the results. Nevertheless, because the state-of-the-art NLP technology is imperfect for named entity and events extractions, as well as for the relationship extractions, we have developed a JavaScript-based online editor to correct errors. The results are stored using the quad RDF (Resource Description Framework) format, which then can be searched via SPARQL, a query language for RDF.

We have experimented with the named-entity extraction of the GATE system using biographical entries on ten Renaissance composers from three different sources: Wikipedia, Oxford Music Online, and the 1911 edition of Grove’s Dictionary of Music and Musicians. The total of 5,441 entities were extracted with the accuracy of 99.24% precision and 98.9% recall. It should be noted, however, that it took over three hours to manually verify and correct the output from the thirty articles; confirming the need for efficient and economical means of correction.

Implications for Musicological Interdisciplinarity

Even though this project concentrates on musicians of the Renaissance, its model can be applied to other time periods and disciplines. Musicians’ lives intersect with artists, writers, clerics, patrons, printers, etc. Combining networks will aid in determining circles of influence and patterns of patronage essential to imagining the past for scholars, teachers, and students of Renaissance culture. As more historical documents are digitized and as the automatic natural language processing improves, a wealth of information that was available but extremely difficult to extract can be more easily retrieved. In fact, this is a pilot project for a more ambitious Human History Project, which aims to create a database for all documented human history, finding references to every individual in the past and able to query any information about them and their relationships with others or to put it another way: “creating a Facebook of the past.”

¹ <http://www.londonlives.org/static/Project.jsp>

Digital Prosopography of Renaissance Musicians

Discovery of Social and Professional Network

Susan Forscher Weiss

Musicology and German & Romance Languages and Literature
The Johns Hopkins University
Robert Lehman Visiting Professor
Harvard Center for Italian Renaissance Studies, Villa I Tatti
sweiss@jh.edu

Ichiro Fujinaga

Music Technology Area, Department of Music Research
Schulich School of Music, McGill University
Centre for Interdisciplinary Research in Music Media and Technology
ich@music.mcgill.ca

Goals

- ◆ To create a social and professional network of Renaissance musicians
- ◆ To create a database to be able to study the network
- ◆ To make connections and discover relationships or answer questions related to dates, geographies, professions, etc.

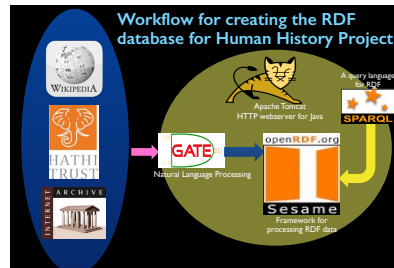
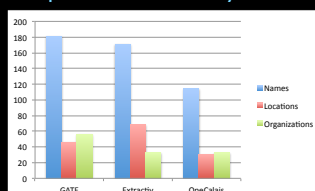
Example queries: Not easily answered by Google, Wikipedia, and other traditional methods

- ◆ Which music printers were in business in 1481 in Florence?
- ◆ Which composers were residing in Florence in 1481?
- ◆ Which composers visited Florence in 1481?
- ◆ Which trumpeters were active in Florence between 1481–86?
- ◆ What events took place in 1481 in Florence that required musical performance?
- ◆ What pieces requiring trumpets were performed in Florence in December 1481?

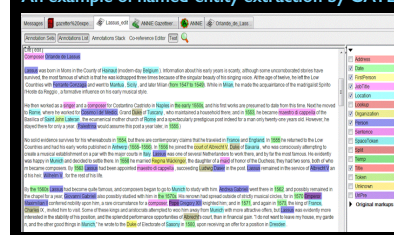
Solutions: Three major tools

- ◆ Natural Language Processing (NLP)
- ◆ Named-entity extraction
- ◆ Events extraction
- ◆ Linked Open Data
- ◆ RDF (Resource Description Framework): a data model
- ◆ SPARQL: a query language for RDF
- ◆ Crowdsourcing

Comparison of Named-Entity Extractors



An example of named-entity extraction by GATE



THE PEABODY INSTITUTE
OF THE JOHNS HOPKINS UNIVERSITY

NATIONAL ENDOWMENT FOR THE HUMANITIES
OFFICE OF DIGITAL HUMANITIES

Named-entity extraction: Experiment I

- ◆ Four Renaissance composers' entries in Wikipedia
- ◆ Using the default Gazetteer (a dictionary)
- ◆ 90.25% precision (the extracted entities were correctly identified)
- ◆ 65.33% recall (the entities in the document that were found)
- ◆ Editing the Gazetteer (~15 minutes / article)
- ◆ 98.39% precision; 91.86% recall
- ◆ Fixing the problem with plurals with the Morphological Analyzer
- ◆ 98.45% precision; 98.45% recall

Named-entity extraction: Experiment II

- ◆ Ten composers x three sources (30) minus the four Wikipedia articles from Experiment I: 26 articles
- ◆ Using the modified Gazetteer and Morphological Analyzer from Experiment I
- ◆ 99.24% precision; 98.9% recall
- ◆ Correction time of the 5,441 entities extracted
- ◆ Average of 3 sec./entity to correct
- ◆ 240 min. to correct 26 articles or about 10 min./ article

Next steps

- ◆ Extract relationships between named entities (e.g., REEL)
- ◆ Create web interface to correct relationships
- ◆ Create web interface to query the network of relationships

Acknowledgements

We would like to acknowledge the wonderful work by Christopher Buchanan on this project. Other team members include: Haoyuan Ji, Peng Jianxiang, Ankit Sarwal, and Chengyuan Zhang. This research has been supported by the NEH Digital Humanities Start-up Grant and the Social Sciences and Humanities Research Council of Canada.

CIR
MMT Centre for Interdisciplinary Research
in Music Media and Technology

McGill

DDMAL

DISTRIBUTED DIGITAL MUSIC
ARCHIVES & LIBRARIES LAB



Social Sciences and Humanities
Research Council of Canada

Conseil de recherches en
sciences humaines du Canada

Canada

Our poster at the AMS/SMT 2014



Activities

Distinguished Lectures

live@CIRMMT
Performance Series

Workshops and
Research Meetings

Research Workshops

Past Workshops

**Workshop on
Human History
Project #1: Digital
prosopography and
linked open data**

Training Workshops

Research Meetings

Seminar Series

Student Colloquia

General Assembly and
Student Symposium

Special Events

Newsletter

Workshop on Human History Project #1: Digital Prosopography and Linked Open Data

— filed under: [Research Workshop](#)

This workshop is organized by CIRMMT Research Axis 3 (Musical information retrieval, archiving and analysis). It will take place on June 21st, in A832 (New Music Building). This workshop is free and open to all. Registration is required.

Registration

Registration is mandatory as seating is limited (35 seats): [Workshop on Human History Project #1- registration](#)

Description

This is the inaugural workshop for initiating a large project called Human History Project, which aims to build a distributed international database of documented human history. Our keynote speaker will be Yves Raimond, who is a pioneer in the Linked Open Data research, especially in the field of music. We will also have presentations by researchers at McGill University who have been working on how to extract, using Natural Language Processing tools and model historical data. We hope to have brainstorming discussions following the presentations.

What	■ Research Workshop
When	Jun 21, 2013 from 10:00 AM to 01:00 PM
Where	A832, New Music Building, 527 Sherbrooke St. West.
Add event to calendar	vCal iCal

Guests / Speakers

- Yves Raimond, BBC, UK
- Ichiro Fujinaga, CIRMMT, Schulich School of Music, McGill University
- Michel Gagnon, École Polytechnique de Montréal
- Yu Hua, Computing Science, McGill University
- Matt Milner, McGill Digital Humanities, McGill University
- Jin Xing, Computing Science, McGill University

Keynote presentation by Yves Raimond

LINKED DATA, ARCHIVES AND THE BBC

Workshop #1 website (2013)



Activities

[Distinguished Lectures](#)

[live@CIRMMT](#)

[Performance Series](#)

[Workshops and
Research Meetings](#)

[Research Workshops](#)

[Past Workshops](#)

**The second
international
workshop on
Human History
Project: Natural
language
processing and big
data**

[Training Workshops](#)

[Research Meetings](#)

[Seminar Series](#)

[Student Colloquia](#)

[General Assembly and
Student Symposium](#)

The second international workshop on Human History Project: Natural language processing and big data

— filed under: [Research Workshop](#)

This workshop is organized by Research Axis 2 [Music information research].

Registration

Registration: Registration for the second international workshop on Human History Project

Description

This is the second workshop on a large project called Human History Project, which aims to build a distributed international database of documented human history using Natural Language Processing tools and Linked Open Data to model historical data.

Guests

Jason Boyd, Department of English, Ryerson University

Serge ter Braake, Faculties of Humanities, University of Amsterdam

Ichiro Fujinaga, CIRMMT, Schulich School of Music, McGill University

Sergio Oramas, Music Technology Group, Universitat Pompeu Fabra

Matthew Milner, Department of History, McGill University

Susan Forscher Weiss, Peabody Conservatory of Music, Johns Hopkins University

What	■ Research Workshop
When	Dec 12, 2015 from 09:00 AM to 02:30 PM
Where	A832, Elizabeth Wirth Music Building, 527 Sherbrooke St. West.
Add event to calendar	vCal iCal

Workshop #2 website (2015)

Other Online Biographies

- Trove (Australia)
 - Digitized newspaper articles with "crowd sourcing"
 - 154,000 names
 - Relationships
- Orlando Project: Women's Writing in the British Isles
 - 1,200 names
- Ming Qing Women's Writings (1368–1911)
 - 6,000 names
- Buddhist Authority Database
 - 23,000 names

Digital Prosopography Projects

- Prosopography of Byzantine World (King's College: 1998–)
 - Originally on CD-ROM (2001)
- Prosopography of Anglo-Saxon England (PASE) (King's College and Cambridge: 2000–)
 - More than 30,000 names
- Clergy of the Church of England Database (1540–1835) (King's College: 1999–)
 - Berkeley Prosopography Services (2009–)
 - Hellenistic Babylonia cuneiforms

The Prosopography of Renaissance Singers

Experiments

- Name extraction from Baker's and other music biographical dictionaries
- Selecting Renaissance musicians
- Named-entity extraction of 10 Renaissance composers
 - Wikipedia
 - Oxford Music Online
 - Grove's Dictionary (1911 edition)

GATE: Named-entity extraction

Experiment with 10 Composers

Named-entity extraction from three sources

Processed articles from Wikipedia, Oxford Music Online, and the 1911 edition of Grove's for each of the ten composers

The results

- Named-entity extraction with 30 articles
 - 10 composers x 3 sources
 - Result: 99.24% precision; 98.9% recall
- Correction time
 - 5,441 entities extracted
 - Average of 3 sec./entity to check and correct
 - 4 hours to correct 30 articles or about 8 min./article

Where is all this data headed? Human History Project

- Create RDF (Resource Description Framework) database
- To correct errors, use crowd sourcing (or more specifically expert sourcing or "grad sourcing")
- Use SPARQL for query
 - SPARQL: a query language for RDF

Workflow for creating the RDF database for Human History Project

Natural Language Processing: Relation Extraction

Heinrich Glaser (also Glaserian) (June 1488 – 28 March 1563) was a Swiss music theorist, poet and humanist. He was born in Mollis (in the canton of Glarus, hence his name) and died in Freiburg.

After a thorough early training in music, he enrolled in the University of Cologne, where he studied theology, philosophy, and mathematics as well as music. It was there that he wrote a famous poem as a tribute to Emperor Maximilian I. Shortly afterwards, in Basel, he met Erasmus and the two humanists became lifelong friends.

Web Editor for Crowd Sourcing

SPARQL Query

Who was born in the 15th century and died in Freiburg?

```

SELECT ?Person ?Birth_date WHERE {
  ?Person <http://www.ontologydesignpatterns.org/ont/foaf/01/name> "Heinrich Glaser" .
  ?Person <http://www.ontologydesignpatterns.org/ont/foaf/01/birthdate> ?Birth_date .
  FILTER (
    ?Birth_date > "1480-01-01"^^xsd:date &&
    ?Birth_date < "1560-01-01"^^xsd:date
  )
}
```

Summary

- Human History Project (HHP)
- Digital Prosopography of Renaissance Musicians
- Named-entity extraction
- Event / Relation extraction
- Prototype for crowd-sourced correction web interface
- Need for simpler query interface for SPARQL

Acknowledgements

Christopher Bailison
Neil Fenley
Andrew Harkinson
Yu Hua
Hayuan Ji
Jin Jin
Sailing Li
Matt Pilsner
Gabriel Vigliani
Jin Xing



Basic Named-Entity Categories

- Entity names: Organisation, person, location, place names, (dates, events, resources)
- Temporal expressions: Date, time
- Number expressions: Money, percent, age, weight, distance

Experiments with GATE

GATE *general architecture for text engineering*

- Named-entity extraction with a total of ten Renaissance composers' entries in Wikipedia, Oxford Music Online, and Grove's (1911)
- Two experiments
 - Experiment 1: Optimizing the configuration (four composers, Wikipedia articles only)
 - Experiment 2: Calculate accuracy and correction time (ten composers and three sources)

Experiment 1

Processed Wikipedia articles of four composers

Galilei Luzzini Scarlatti Vivaldi

The results of experiment 1

- Named-entity extraction with four Renaissance composers' entries in Wikipedia
- Using the default Gazetteer (a dictionary)
 - 98.3% precision (the extracted entities were correctly identified)
 - 45.3% recall (the entities in the document that were found)
- Editing the Gazetteer (~15 minutes / article)
 - 98.3% precision / 91.8% recall
- Problem with plurals
 - Fixed with the Morphological Analyzer (finding root form of words)
 - 98.4% precision / 98.4% recall

Experiment 2

Galilei Luzzini Scarlatti Vivaldi

Baker's A Biographical Dictionary of Musicians (1900) from the Internet Archive

Baker's A Biographical Dictionary of Musicians (1900)

Isidoro (Isidoro, Isidoro, Isidoro, 1490, 1500)
Isidoro (Isidoro, Isidoro, Isidoro, 1490, 1500)
Isidoro (Isidoro, Isidoro, Isidoro, 1490, 1500)
Isidoro (Isidoro, Isidoro, Isidoro, 1490, 1500)
Isidoro (Isidoro, Isidoro, Isidoro, 1490, 1500)
Isidoro (Isidoro, Isidoro, Isidoro, 1490, 1500)
Isidoro (Isidoro, Isidoro, Isidoro, 1490, 1500)
Isidoro (Isidoro, Isidoro, Isidoro, 1490, 1500)
Isidoro (Isidoro, Isidoro, Isidoro, 1490, 1500)
Isidoro (Isidoro, Isidoro, Isidoro, 1490, 1500)

Names extracted from other biographical dictionaries

- Baker's Dictionary of Musicians, 1914 ed.
- 1449 names
- British Musical Biography, 1897 ed.
- 1476 names
- Biographical Dictionary of Musicians (Cummings), 1892 ed.
- 1,711 names
- 136 Renaissance musicians
- Grove's Dictionary of Music and Musicians, 3 vols, 1879 ed.
- 1,773 names

The final results of the previous experiments and the progress

- Total number of Renaissance musicians
 - 304 (Baker) + 136 (Cummings) = 340 (349 different names)
- Progress since the last report
 - Tested various named-entity extraction techniques
 - Using Natural Language Processing (NLP) software
 - GATE, Alchemy, OpenCalais, etc.

Results and current research

- Total number of Renaissance musicians
 - 304 (Baker) + 136 (Cummings) = 340 (349 different names)
 - Use these names as seeds to find other musicians' entries
- Testing various Named-Entity Extraction techniques
- Several software available
 - OpenCalais, GATE, Alchemy, etc.

Alchemy (a screenshot)

Comparison of Named-Entity Extraction Performance

Tool	Precision	Recall	F1 Score
GATE	~98%	~45%	~65%
Alchemy	~98%	~92%	~95%
OpenCalais	~98%	~92%	~95%

Next step: Event extractions

- Event ontologies
 - Event
 - LODE (Linking Open Description of Events)
- Quad RDF representation
- The importance of provenance

Event

LODE

Linking Open Description of Events

"... an ontology for publishing descriptions of historical events as Linked Data, and for mapping between other event-related vocabularies and ontologies."

Term Name	Type	Definition
event	class	"Something that happened." An event is something that is a time instant or is supported by a historical source.
place	property	A named or relatively specified place that is where an event happened.
when	property	An historical instant or interval of time that is when an event happened.
about	property	An interval of time that can be precisely described using calendar dates and clock times.
source	property	An external source of events that is a historical document or report that is where an event happened.
related	property	An external region of events that is a geographical point or region that is where an event happened.
related	property	An external, whole, or temporal object involved in an event.
related	property	An event involved in an event.

Representing events in RDF

- Statement 123: Anna Maria Mozart gives birth to a son
- Statement 124: Statement 123 occurred in 1756
- Statement 125: Statement 124 was written in Wikipedia
- Statement 126: Statement 124 was written on 2010/11/19
- Statement 127: Statement 123 was written in a book
- Statement 128: Statement 123 took place in Salzburg

Provenance Ontologies

- Open Provenance Model
- Provenir Ontology
- Provenance Vocabulary
- etc.
- Provenance Vocabulary Mapping (W3C)

Other tools to try

- Apache OpenNLP
- Stanford NLP Tools
- AFNER: Named Entity Recognition (Australia)
- FRED: "a tool for automatically producing RDF/OWL ontologies and linked data from natural language sentences" (Italy)
- Apache Stanbol: "... provides a set of reusable components for semantic content management"

Existing "structured" data sources

- Freebase / Knowledge Graph
- DBpedia
- Yago (Max Planck Institute)
- Need to add provenances
- VIAF (Virtual International Authority Files): unique ID/URI
 - J.S. Bach URI: (<http://viaf.org/viaf/12304462/>)
- WorldCat Identities (30 million names)
 - 100 billion total (10¹¹)
- Getty Thesaurus of Geographic Names (Montréal)

Event extraction applications

- Personalized news systems
- Risk analysis
- Decision-making support tools
- Biological event detectors from academic papers
- Algorithmic stock trading

Event extraction research

- "On-line news event and tracking" (1998)
- "A system for news event detection" (2003)
- "Real-time news event extraction for global crises monitoring" (2008)
- "An overview of event extraction from text" (2011).
Workshop on Detection, Representation, and Exploitation of Events in the Semantic Web (DeRIVE 2011)
- "A survey of techniques for event detection in Twitter" (2013)

Workshop on Detection, Representation, and Exploitation of Events in the Semantic Web

- Started in 2011
- "The goal of this workshop is to bring together those different areas in the recent surge of research on the use of events as a key concept for representing and organizing knowledge on the Web." (<http://derive2013.wordpress.com>)

HHP: Multidisciplinary Research

- History
- Linguistics
- Computer Science
- Information Science
- Geography
- Philosophy ("what sort of things are events?")
- Onomastics: Study of proper names
- Toponymy: Study of place names
- Hydronymy: Study of a proper name of a body of water

HHP: Future Plans

- Suggestion for workshops
- Build an infrastructure at McGill?
- Grant applications
- Iterative improvements over the years (next 50-100 years)
- Logo design / designers?

Basic Entity Categories

- Entity names (ENAMEX): Organization, person, location, place names (lake, river, mountain)
- Temporal expression (TIMEX): Date, time
- Number expressions (NUMEX): Money, percent, age, weight, distance

Named-entity Recognition

- The term "named-entity recognition (NER)" first used in a paper by Ralph Grishman (Computer Science, NYU). Proceedings of the workshop on Human Language Technology (held at Rambouillet, March 8-11, 1994)
- NER was developed as part of the evaluation for Message Understanding Conference (MUC) which began in 1987. As part of the evaluation suite MUC-6 using TREC corpora
- MUC (began in 1987) have been organized by NRAD, the Research Development, Testing and Evaluation division of the Naval Command Control and Ocean Surveillance Center with the support of DARPA (the Defense Advanced Research Projects Agency)
- Message Understanding Conferences were forum "to assess and to foster research on the automated analysis of military messages..." (Grishman and Sundheim 1996, 466)

History of Message Understanding Conferences (MUC)

- "Although called 'conferences', the distinguishing characteristic of the MUCs are not the conferences themselves, but the evaluations to which participants must submit in order to be permitted to attend the conference." (Grishman and Sundheim 1996, 466)
- "MUC-1 (1987) was basically exploratory.... By MUC-2 (1989), the task" was to "fill a template with information about the event, such as the type of event, the agent, the time and place, the effect, etc."

